# Building the Next Generation of Parallel Applications

Michael A. Heroux

Scalable Algorithms Department

Sandia National Laboratories, USA

Sandia National Laboratories

# A Brief Personal Computing History

**1988 - 1997**

**1993 - 2008**

```
CMIC$ DO ALL VECTOR IF (N .GT. 800)
CMIC$1   SHARED(BETA, N, Y, Z)
CMIC$2   PRIVATE(I)
CDIR$ IVDEP
      do 15 i = 1, n
        z(i) = beta * y(i)
  15    continue
     endif
```

```
#include <mpi.h>
int main(int argc, char *argv[]) {
// Initialize MPI
  MPI_Init(&argc,&argv);
  int rank, size;
  MPI_Comm_rank(MPI_COMM_WORLD, &rank);
  MPI_Comm_size(MPI_COMM_WORLD, &size);
```

Sandia National Laboratories

# 2008 - Present

**Unification and composition:**
**-Vectorization**
**-Threading**
**- Multiprocessing**

```
#include <mpi.h>
#include <omp.h>
int main(int argc, char *argv[]) {
// Initialize MPI
 MPI_Init(&argc,&argv);
 int rank, size;
 MPI_Comm_rank(MPI_COMM_WORLD, &rank);
 MPI_Comm_size(MPI_COMM_WORLD, &size);
…
#pragma omp parallel
{
    double localasum = 0.0;
#pragma omp for
    for (int j=0; j< MyLength_; j++) localasum += std::abs(from[j]);
#pragma omp critical
    asum += localasum;
}
```

```
#include <thrust/host_vector.h>
#include <thrust/device_vector.h>

thrust::device_vector<int> vd(10, 1);
thrust::host_vector<int> vh(10,1);
```
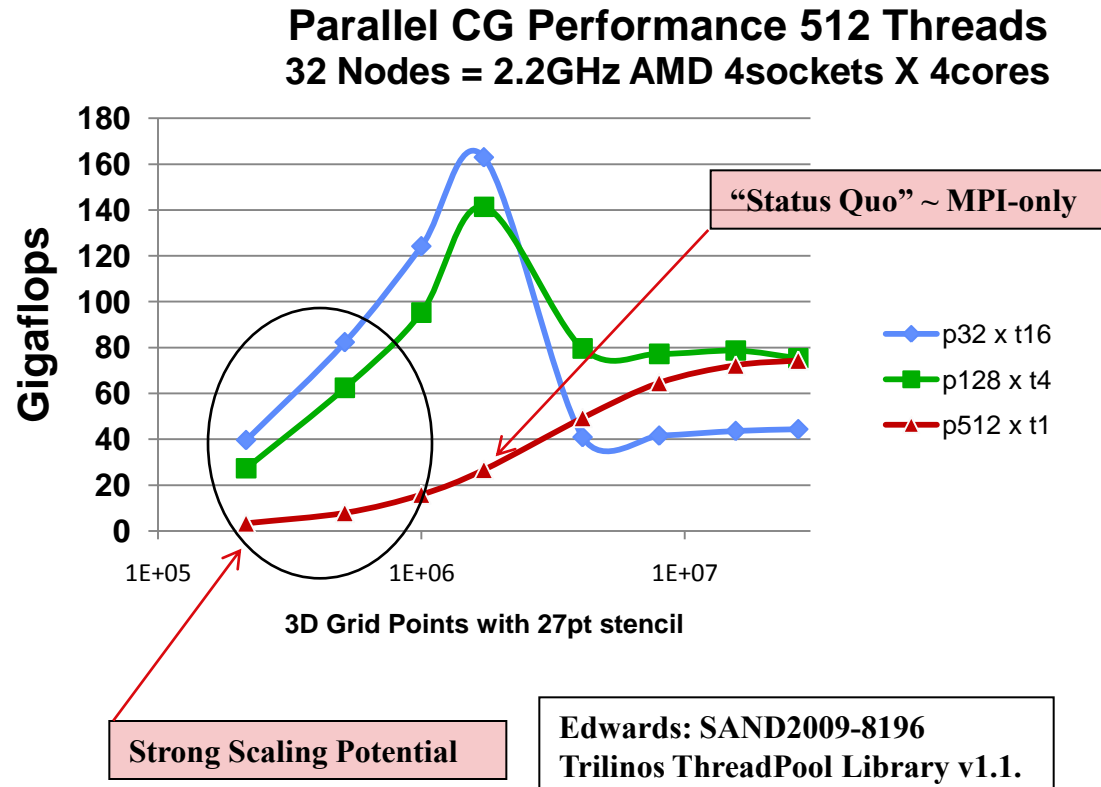
Sandia National Laboratories

# Quiz (True or False)

1. MPI-only has the best parallel performance.
2. Future parallel applications will not have MPI_Init().
3. All future programmers will need to write parallel code.
4. Use of "markup", e.g., OpenMP pragmas, is the least intrusive approach to parallelizing a code.
5. DRY is not possible across CPUs and GPUs
6. GPUs are a harbinger of CPU things to come.
7. Checkpoint/Restart will be sufficient for scalable resilience.
8. Resilience will be built into algorithms.
9. MPI-only and MPI+X can coexist in the same application.
10. Kernels will be different in the future.

# Basic Exascale Concerns: Trends, Manycore

- Stein's Law: *If a trend cannot continue, it will stop.*

  Herbert Stein, chairman of the Council of Economic Advisers under Nixon and Ford.

- Trends at risk:
  - Power.
  - Single core performance.
  - Node count.
  - Memory size & BW.
  - Concurrency expression in existing Programming Models.

**Parallel CG Performance 512 Threads
32 Nodes = 2.2GHz AMD 4sockets X 4cores**

**"Status Quo" ~ MPI-only**

Legend:
- p32 x t16
- p128 x t4
- p512 x t1

Y-axis: **Gigaflops** (0, 20, 40, 60, 80, 100, 120, 140, 160, 180)

X-axis: 1E+05, 1E+06, 1E+07 — **3D Grid Points with 27pt stencil**

**Strong Scaling Potential**

**Edwards: SAND2009-8196
Trilinos ThreadPool Library v1.1.**

**One outcome: Greatly increased interest in OpenMP**

Sandia National Laboratories

# Implications

- MPI-Only is not sufficient, except … much of the time.
- Near-to-medium term:
  - MPI+[OMP|TBB|Pthreads|CUDA|OCL|MPI]
  - Long term, too?
- Long- term:
  - Something hierarchical, global in scope.
- Conjecture:
  - Data-intensive apps need non-SPDM model.
  - Will develop new programming model/env.
  - Rest of apps will adopt over time.
  - Time span: 20 years.

# What Can we Do Right Now?

- Study why MPI was successful.
- Study new parallel landscape.
- Try to cultivate an approach similar to MPI.

# *MPI Impresssions*

# MPI: It Hurts So Good

- ## Observations
  - "assembly language" of parall...
  - lowest common denomi...
    - portable across ar...
  - upfront effort r...
    - system
    - e...
- C...

So What Would Life Be Like Without MPI?

Looking Forward to a New Age of Large-Scale Parallel Programming and the Demise of MPI
...hopes and dreams of an HPC educator

$$F(n) = \begin{cases} 0 & n=0 \\ 1 & n=1 \\ F(n-1)+F(n-2) & n>1 \end{cases}$$

```
Serial C
long fib_serial(long n)
{   if (n < 2) return n;
    return fib_serial(n-1) + fib_serial(n-2);
}
```

```
OpenMP 3.0
long fib_parallel(long n)
{   long x, y;
    if (n < 2) return n;
    #pragma omp task default(none) shared(x,n)
    {   x = fib_parallel(n-1);
    }   y = fib_parallel(n-2);
    #pragma omp taskwait
    return (x+y);
}
```

```
Cilk++
long fib_parallel(long n)
{   long x, y;
    if (n < 2) return n;
    x = cilk_spawn fib_parallel(n-1);
    y = fib_parallel(n-2);
    cilk_sync;
    return (x+y);
}
```

```
def fib_ser
    var x,y; n.
    if (n < 2) th
    cobegin {
        x=fib_serial(
        y=fib_serial(
    }
    return x+y;
}
```

```
fib_parallel(n: ZZ): ZZ requ
    if n < 2 then n else fib_pa
```

RELAX. WE'RE TOO BIG TO FAIL.

Fortran
MPI
C

Brad Chamberlain, Cray, PPOPP'06, http://chapel.cray.com/publications/ppopp06-slides.pdf

# *MPI Reality*

# dft_fill_wjdc.c



## Tramonto WJDC Functional

- New functional.
- Bonded systems.
- 552 lines C code.

**WJDC-DFT (Werthim, Jain, Dominik, and Chapman) theory for bonded systems.** *(S. Jain, A. Dominik, and W.G. Chapman. Modified interfacial statistical associating fluid theory: A perturbation density functional theory for inhomogeneous complex fluids. J. Chem. Phys., 127:244904, 2007.)* **Models stoichiometry constraints inherent to bonded systems.**

**How much MPI-specific code?**

dft_fill_wjdc.c
MPI-specific
code

**source_pp_g.f**

# MFIX
## Source term for pressure correction

- MPI-callable, OpenMP-enabled.
- 340 Fortran lines.
- No MPI-specific code.
- Ubiquitous OpenMP markup (red regions).

**MFIX: Multiphase Flows with Interphase eXchanges (https://www.mfix.org/)**

# Reasons for MPI Success?

- Portability?                Yes.
- Standardized?             Yes.
- Momentum?               Yes.
- Separation of many Parallel & Algorithms concerns?                Big Yes.


- Once framework in place:
  - Sophisticated physics added as serial code.
  - Ratio of science experts vs. parallel experts: 10:1.


- Key goal for new parallel apps: Preserve this ratio

# Computational Domain Expert Writing MPI Code

# Computational Domain Expert Writing Future Parallel Code

# *Evolving Parallel Programming Model*

# Parallel Programming Model: Multi-level/Multi-device

**network of computational nodes**

| Inter-node/**inter-device** (distributed) parallelism and resource management | ⇐ **Message Passing** |

↓

**Node-local control flow (serial)**

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

**computational node with manycore CPUs and / or GPGPU**

**Intra-node (manycore) parallelism and resource management** ⇐ **Threading**

↓

**Stateless computational kernels run on each core** ⇐ **stateless kernels**

*Adapted from slide of H. Carter Edwards*

Sandia National Laboratories

# Domain Scientist's Parallel Palette

- MPI-only (SPMD) apps:
  - Single parallel construct.
  - Simultaneous execution.
  - Parallelism of even the messiest serial code.
- Next-generation applications:
  - Internode:
    - MPI, yes, or something like it.
    - Composed with intranode.
  - Intranode:
    - Much richer palette.
    - More care required from programmer.
- What are the constructs in our new palette?

Sandia National Laboratories

# Obvious Constructs/Concerns

- Parallel for:
  - No loop-carried dependence.
  - Rich loops.
- Parallel reduce:
  - Couple with other computations.
  - Concern for reproducibility.

# Other construct: Pipeline

- Sequence of filters.
- Each filter is:
  - Sequential (grab element ID, enter global assembly) or
  - Parallel (fill element stiffness matrix).
- Filters executed in sequence.
- Programmer's concern:
  - Determine (conceptually): Can filter execute in parallel?
  - Write filter (serial code).
  - Register it with the pipeline.
- Extensible:
  - New physics feature.
  - New filter added to pipeline.

# Other construct: Thread team

- Multiple threads.
- Fast barrier.
- Shared, fast access memory pool.
- Example: Nvidia SM
- X86 more vague, emerging more clearly in future.

# Finite Elements/Volumes/Differences and parallel node constructs

- Parallel for, reduce, pipeline:
  - Sufficient for vast majority of node level computation.
  - Supports:
    - Complex modeling expression.
    - Vanilla parallelism.
- Thread team:
  - Complicated.
  - Requires true parallel algorithm knowledge.
  - Useful in solvers.

Sandia National Laboratories

# Preconditioners for Scalable Multicore Systems

## Charon Timing Breakdown on TLCC
### Strong Scaling 28M Unknowns



Legend:
- Charon minus solver
- Solve time due to iter increase
- Solve time due to iter cost
- Preconditioner setup

x-axis: # Procs (128, 256, 512, 1024, 2048, 4096)
y-axis: Percent Time (0% – 100%)

**Strong scaling of Charon on TLCC (P. Lin, J. Shadid 2009)**

## # Linear Solver Iterations per Newton Step



Values: 111, 117, 117, 125, 129, 153

# MPI Ranks: 128, 256, 512, 1024, 2048, 4096

- Observe: Iteration count increases with number of subdomains.
- With scalable threaded triangular solves
  - Solve triangular system on larger subdomains.
  - Reduce number of subdomains.
- Goal:
  - Better kernel scaling (threads vs. MPI processes).
  - Better convergence, More robust.
- Note: App (-solver) scales very well in MPI-only mode.
- Exascale Potential: Tiled, pipelined implementation.

| MPI Tasks | Threads | Iterations |
|---|---|---|
| 4096 | 1 | 153 |
| 2048 | 2 | 129 |
| 1024 | 4 | 125 |
| 512 | 8 | 117 |
| 256 | 16 | 117 |
| 128 | 32 | 111 |

*Factors Impacting Performance of Multithreaded Sparse Triangular Solve*, Michael M. Wolf and Michael A. Heroux and Erik G. Boman, VECPAR 2010, to appear.

Sandia National Laboratories

# Level Set Triangular Solver



**L**

**DAG**

Triangular Solve:
- Critical Kernel
  - MG Smoothers
  - Incomplete IC/ILU
- Naturally Sequential
- Building on classic algorithms:
  - Level Sched:
    - circa 1990.
    - Vectorization.
  - New: Generalized.

$$\tilde{L} = PLP^T = \begin{bmatrix} D_1 & & & & \\ A_{2,1} & D_2 & & & \\ A_{3,1} & A_{3,2} & D_3 & & \\ \vdots & \vdots & \vdots & \ddots & \\ A_{l,1} & A_{l,2} & A_{l,3} & \dots & D_l \end{bmatrix}$$

Permuted System

$$\tilde{x}_1 = D_1^{-1}\tilde{y}_1$$

$$\tilde{x}_2 = D_2^{-1}\left(\tilde{y}_2 - A_{2,1}\tilde{x}_1\right)$$

Multi-step Algorithm

$$\vdots \quad \vdots \qquad \vdots$$

$$\tilde{x}_l = D_l^{-1}\left(\tilde{y}_l - A_{l,1}\tilde{x}_1 - \dots - A_{l,l-1}\tilde{x}_{l-1}\right)$$

Sandia National Laboratories

# Triangular Solve Results

| Name | N | nnz | N/nlevels | application area |
|---|---|---|---|---|
| asic680ks | 682,712 | 2,329,176 | 13932.9 | circuit simulation |
| cage12 | 130,228 | 2,032,536 | 1973.2 | DNA electrophoresis |
| pkustk04 | 55,590 | 4,218,660 | 149.4 | structural engineering |
| bcsstk32 | 44,609 | 2,014,701 | 15.1 | structural engineering |

**Passive (PB) vs. Active (AB) Barriers: Critical for Performance**



Nehalem

Istanbul

**AB + No Thread Affinity (NTA) vs. AB + Thread Affinity (TA) : Also Helpful**

Level sets: Trilinos/Isorropia          Core Kernel Timings: Trilinos/Kokkos.

Sandia National Laboratories

# Thread Team Advantanges

- **Qualitatively better algorithm:**
  - **Threaded triangular solve scales.**
  - **Fewer MPI ranks means fewer iterations, better robustness.**
- **Exploits:**
  - **Shared data.**
  - **Fast barrier.**
  - **Data-driven parallelism.**

# *Placement and Migration*

Sandia National Laboratories

# Placement and Migration

- MPI:
  - Data/work placement clear.
  - Migration explicit.
- Threading:
  - It's a mess (IMHO).
  - Some platforms good.
  - Many not.
  - Default is bad (but getting better).
  - Some issues are intrinsic.

Sandia National Laboratories

# Data Placement on NUMA

- Memory Intensive computations: Page placement has huge impact.

- Most systems: First touch (except LWKs).

- Application data objects:
  - Phase 1: Construction phase, e.g., finite element assembly.
  - Phase 2: Use phase, e.g., linear solve.

- Problem: First touch difficult to control in phase 1.

- Idea: Page migration.
  - Not new: SGI Origin.  Many old papers on topic.

# Data placement experiments

- MiniApp: HPCCG (Mantevo Project)

- Construct sparse linear system, solve with CG.

- Two modes:
  - Data placed by assembly, not migrated for NUMA
  - Data migrated using parallel access pattern of CG.

- Results on dual socket quad-core Nehalem system.

Sandia
National
Laboratories

# Weak Scaling Problem



**Weak Scaling**
**Dim 260K Per core**

- MPI and conditioned data approach comparable.
- Non-conditioned very poor scaling.

# Page Placement summary

- MPI+OpenMP (or any threading approach) is best overall.

- But:
  - Data placement is big issue.
  - Hard to control.
  - Insufficient runtime support.

- Current work:
  - Migrate on next-touch (MONT).
  - Considered in OpenMP (next version).
  - Also being studied in Kitten (Kevin Pedretti).

- Note: This phenomenon especially damaging to OpenMP common usage.

Sandia National Laboratories

# *Transition: MPI-only to MPI+[X|Y|Z]*

Sandia National Laboratories

# Parallel Machine Block Diagram



- Parallel machine with $p = m * n$ processors:
  - $m$ = number of nodes.
  - $n$ = number of shared memory processors per node.
- Two ways to program:
  - Way 1: $p$ MPI processes.
  - Way 2: $m$ MPI processes with $n$ threads per MPI process.
- New third way:
  - "Way 1" in some parts of the execution (the app).
  - "Way 2" in others (the solver).

Sandia National Laboratories

# Multicore Scaling: App vs. Solver



Tramonto vs. Solver Time on Niagara2: 4-48 Threads



Charon vs Solver Time: 1-16 Cores

## Application:
- Scales well (sometimes superlinear)
- MPI-only sufficient.

## Solver:
- Scales more poorly.
- Memory system-limited.
- MPI+threads can help.

* Charon Results:
  Lin & Shadid TLCC Report

Sandia National Laboratories

# MPI-Only + MPI/Threading: *Ax=b*

**App** Rank 0    **App** Rank 1    **App** Rank 2    **App** Rank 3

App passes matrix and vector values to library data classes

**Lib** Rank 0    **Lib** Rank 1    **Lib** Rank 2    **Lib** Rank 3

All ranks store *A, x, b* data in memory visible to rank 0

**Mem** Rank 0    **Mem** Rank 1    **Mem** Rank 2    **Mem** Rank 3

**Multicore: "PNAS" Layout**

Library solves *Ax=b* using shared memory algorithms on the node.

**Lib** Rank 0

Thread 0   Thread 1   Thread 2   Thread 3

Sandia National Laboratories

# MPI Shared Memory Allocation

Idea:

- Shared memory alloc/free functions:
  - MPI_Comm_alloc_mem
  - MPI_Comm_free_mem
- Predefined communicators:
  MPI_COMM_NODE – ranks on node
  MPI_COMM_SOCKET – UMA ranks
  MPI_COMM_NETWORK – inter node
- Status:
  - Available in current development branch of OpenMPI.
  - First "Hello World" Program works.
  - Incorporation into standard still not certain. Need to build case.
  - Next Step: Demonstrate usage with threaded triangular solve.
- Exascale potential:
  - Incremental path to MPI+X.
  - Dial-able SMP scope.

```
int n = …;
double* values;
 MPI_Comm_alloc_mem(
              MPI_COMM_NODE,  // comm (SOCKET works too)
              n*sizeof(double),      // size in bytes
              MPI_INFO_NULL,    // placeholder for now
              &values);              // Pointer to shared array (out)


// At this point:
// - All ranks on a node/socket have pointer to a shared buffer (values).
// - Can continue in MPI mode (using shared memory algorithms) or
// - Can quiet all but one:
int rank;
MPI_Comm_rank(MPI_COMM_NODE, &rank);
if (rank==0) { // Start threaded code segment, only on rank 0 of the node
…
}


 MPI_Comm_free_mem(MPI_COMM_NODE, values);
```

**Collaborators: B. Barrett, Brightwell, Wolf - SNL; Vallee, Koenig - ORNL**

Sandia National Laboratories

# *Resilient Algorithms*

# My Luxury in Life (wrt FT/Resilience)

The privilege to think of a computer as a *reliable, digital* machine.

"At 8 nm process technology, it will be harder to tell a 1 from a 0."

(W. Camp 2008, 2010)

Sandia National Laboratories

# Users' View of the System Now

- "All nodes up and running."
- Certainly nodes fail, but invisible to user.
- No need for me to be concerned.
- Someone else's problem.

Sandia
National
Laboratories

- Nodes in one of four states.

1. Dead.
2. Dying (perhaps producing faulty results).
3. Reviving.
4. Running properly:
   a) Fully reliable or…
   b) Maybe still producing an occasional bad result.

# Faults: Hard vs. Soft

- Hard:
  - Program flow interrupted.
  - Majority of faults.
  - Presently handled by (global) checkpoint/restart.
  - Numerous papers on alternatives.
- Soft:
  - Program flow continues.
  - Minor perturbations in data state:
    - Incorrect address lookup (but still in user scope).
    - Incorrect FP value.

Sandia National Laboratories

# Algorithm-Based (Hard) Fault Tolerance

- Numerous approaches.
- Most common strategies:
  - Meta data:
    - Embed meta data into user-defined data structures.
    - Manage fault detection, recovery manually.
  - Algorithm results validation:
    - Use known algorithm properties.
    - Validate computed to known (e.g., residual check).
- Note: A lack of app awareness.

Sandia National Laboratories

Common Approach to FT (Diplomacy Analogy)

# Hard Error Futures

- C/R will continue as dominant approach:
  - Global state to global file system OK for small systems.
  - Large systems: State control will be localized, use SSD.
- Checkpoint-less restart:
  - Requires full vertical HW/SW stack co-operation.
  - Very challenging.
  - Stratified research efforts not effective.

# Soft Error Futures

- Soft error handling: A legitimate algorithms issue.
- Programming model, runtime environment play role.

Sandia National Laboratories

# Consider GMRES as an example of how soft errors affect correctness

- Basic Steps
  1) Compute Krylov subspace (preconditioned sparse matrix-vector multiplies)
  2) Compute orthonormal basis for Krylov subspace (matrix factorization)
  3) Compute vector yielding minimum residual in subspace (linear least squares)
  4) Map to next iterate in the full space
  5) Repeat until residual is sufficiently small

- More examples in Bronevetsky & Supinski, 2008

Sandia National Laboratories

# Why GMRES?

- Many apps are implicit.

- Most popular (nonsymmetric) linear solver is preconditioned GMRES.

- Only small subset of calculations need to be reliable.

  – GMRES is iterative, but also direct.

# Every calculation matters

- Small PDE Problem: Dim 21K, Nz 923K.
- ILUT/GMRES
- Correct computation 35 Iters: 343M FLOPS
- Two examples of a single bad floating point op

| Description | Iterations | FLOPS | Recursive Residual Error | Solution Error |
|---|---|---|---|---|
| All Correct Calcs | 35 | 343M | 4.6e-15 | 1.0e-6 |
| Iter=2, y[1] += 1.0 SpMV incorrect Ortho subspace | 35 | 343M | 6.7e-15 | 3.7e+3 |
| Q[1][1] += 1.0 Non-ortho subspace | N/C | N/A | 7.7e-02 | 5.9e+5 |

Sandia National Laboratories

# One possible approach is transactional computation

- Database transactions: atomic

- Transactional memory: atomic memory operation

- Transactional computation:
  - Designated sensitive computation region (orthogonalization step in GMRES)
  - Guarantee accurate computation or notify user.

Sandia National Laboratories

# Needs to be coupled with SW-enabled guaranteed data regions

- User-designated reliable data region
- Extra protection to improve reliable data storage and transfer
- Examples
  - Original input data (needed for verification)
  - Linear solver: *A, x, b*
  - Orthogonal vectors for GMRES
- OpenMP pragma-enabled?

Sandia National Laboratories

# Goal

- Algorithms well-conditioned wrt soft failure.

- Now:
  - Single soft error produces erroneous results.

- Goal:
  - Correct results always.
  - Cost increase proportional to number of soft errors.

- Note: These are just two approaches to ABFT.

# *Software Development and Delivery*

Sandia National Laboratories

# Compile-time Polymorphism
## Templates and Sanity upon a shifting foundation

Software delivery:

- Essential Activity

How can we:

- Implement mixed precision algorithms?
- Implement generic fine-grain parallelism?
- Support hybrid CPU/GPU computations?
- Support extended precision?
- Explore redundant computations?
- Prepare for both exascale "swim lanes"?

C++ templates only sane way:

- Moving to completely templated Trilinos libraries.
- Other important benefits.
- **A usable stack exists now in Trilinos**.

Template Benefits:
- Compile time polymorphism.
- True generic programming.
- No runtime performance hit.
- Strong typing for mixed precision.
- Support for extended precision.
- Many more…

Template Drawbacks:
- Huge compile-time performance hit:
  - But good use of multicore :)
  - Eliminated for common data types.
- Complex notation:
  - Esp. for Fortran & C programmers).
  - Can insulate to some extent.

Sandia National Laboratories

# Solver Software Stack

**Phase I packages: SPMD, int/double**  **Phase II packages: Templated**

*Trilinos*

| | | Sensitivities (Automatic Differentiation: Sacado) | |
|---|---|---|---|
| **Optimization**<br>Unconstrained:<br>Constrained: | Find $u \in \Re^n$ that minimizes $g(u)$<br>Find $x \in \Re^m$ and $u \in \Re^n$ that minimizes $g(x,u)$ s.t. $f(x,u) = 0$ | | **MOOCHO** |
| **Bifurcation Analysis** | Given nonlinear operator $F(x,u) \in \Re^{n+m}$<br>For $F(x,u) = 0$ find space $u \in U \ni \frac{\partial F}{\partial x}$ | | **LOCA** |
| **Transient Problems**<br>DAEs/ODEs: | Solve $f(\dot{x}(t), x(t), t) = 0$<br>$t \in [0,T], x(0) = x_0, \dot{x}(0) = x_0'$<br>for $x(t) \in \Re^n, t \in [0,T]$ | | **Rythmos** |
| **Nonlinear Problems** | Given nonlinear operator $F(x) \in \Re^m \to \Re$<br>Solve $F(x) = 0$  $x \in \Re^n$ | | **NOX** |
| **Linear Problems**<br>Linear Equations:<br>Eigen Problems: | Given Linear Ops (Matrices) $A, B \in \Re^{m \times n}$<br>Solve $Ax = b$ for $x \in \Re^n$<br>Solve $A\nu = \lambda B\nu$ for (all) $\nu \in \Re^n$, $\lambda \in$ | | **Anasazi**<br>**Ifpack, ML, etc...**<br>**AztecOO** |
| **Distributed Linear Algebra**<br>Matrix/Graph Equations:<br>Vector Problems: | Compute $y = Ax; A = A(G); A \in \Re^{m \times n}, G \in \Im^{m \times n}$<br>Compute $y = \alpha x + \beta w; \alpha = \langle x, y \rangle; x, y \in \Re^n$ | | **Epetra**<br>**Teuchos** |

# Solver Software Stack

| Phase I packages | Phase II packages | Phase III packages: Manycore*, templated |
|---|---|---|

| **Optimization**<br>  **Unconstrained:**<br>  **Constrained:** | Find $u \in \Re^n$ that minimizes $g(u)$<br>Find $x \in \Re^m$ and $u \in \Re^n$ that minimizes $g(x,u)$ s.t. $f(x,u) = 0$ | **Sensitivities**<br>**(Automatic Differentiation: Sacado)** | **MOOCHO** | |
|---|---|---|---|---|
| **Bifurcation Analysis** | Given nonlinear operator $F(x,u) \in \Re^{n+m}$<br>For $F(x,u) = 0$ find space $u \in U \ni \dfrac{\partial F}{\partial x}$ | | **LOCA** | **T-LOCA** |
| **Transient Problems**<br>  **DAEs/ODEs:** | Solve $f(\dot{x}(t), x(t), t) = 0$<br>$t \in [0, T], x(0) = x_0, \dot{x}(0) = x_0'$<br>for $x(t) \in \Re^n, t \in [0, T]$ | | **Rythmos** | |
| **Nonlinear Problems** | Given nonlinear operator $F(x) \in \Re^m \to \Re$<br>Solve $F(x) = 0$ $x \in \Re^n$ | | **NOX** | **T-NOX** |
| **Linear Problems**<br>  **Linear Equations:**<br>  **Eigen Problems:** | Given Linear Ops (Matrices) $A, B \in \Re^{m \times n}$<br>Solve $Ax = b$ for $x \in \Re^n$<br>Solve $A\nu = \lambda B\nu$ for (all) $\nu \in \Re^n$, $\lambda \in$ | | **Anasazi** | |
| | | | **AztecOO**<br>**Ifpack,**<br>**ML, etc...** | **Belos***<br>**T-Ifpack*,**<br>**T-ML*, etc.** |
| **Distributed Linear Algebra**<br>  **Matrix/Graph Equations:**<br>  **Vector Problems:** | Compute $y = Ax; A = A(G); A \in \Re^{m \times n}, G \in \Im^{m \times n}$<br>Compute $y = \alpha x + \beta w; \alpha = \langle x, y \rangle; x, y \in \Re^n$ | | **Epetra** | **Tpetra***<br>**Kokkos** |
| | | | **Teuchos** | |

# *Trilinos/Kokkos Node API*

# Generic Shared Memory Node

- Abstract inter-node comm provides DMP support.
- Need some way to portably handle SMP support.
- Goal: allow code, once written, to be run on any parallel node, regardless of architecture.
- Difficulty #1: Many different memory architectures
  - Node may have multiple, disjoint memory spaces.
  - Optimal performance may require special memory placement.
- Difficulty #2: Kernels must be tailored to architecture
  - Implementation of optimal kernel will vary between archs
  - No universal binary ➔ need for separate compilation paths

Sandia National Laboratories

# Kokkos Node API

- Kokkos provides two main components:
  - Kokkos memory model addresses Difficulty #1
    - Allocation, deallocation and efficient access of memory
    - compute buffer: special memory used for parallel computation
    - New: Local Store Pointer and Buffer with size.
  - Kokkos compute model addresses Difficulty #2
    - Description of kernels for parallel execution on a node
    - Provides stubs for common parallel work constructs
    - Currently, parallel for loop and parallel reduce
- Code is developed around a polymorphic Node object.
- Supporting a new platform requires only the implementation of a new node type.

Sandia National Laboratories

# Kokkos Memory Model

- A generic node model must at least:
  - support the scenario involving distinct device memory
  - allow efficient memory access under traditional scenarios
- Nodes provide the following memory routines:

```
ArrayRCP<T> Node::allocBuffer<T>(size_t sz);
void        Node::copyToBuffer<T>(  T * src,
                                  ArrayRCP<T>  dest);
void        Node::copyFromBuffer<T>(ArrayRCP<T> src,
                                  T * dest);
ArrayRCP<T> Node::viewBuffer<T> (ArrayRCP<T> buff);
void        Node::readyBuffer<T>(ArrayRCP<T> buff);
```

# Kokkos Compute Model

- How to make shared-memory programming generic:
  - Parallel reduction is the intersection of `dot()` and `norm1()`
  - Parallel for loop is the intersection of `axpy()` and mat-vec
  - We need a way of fusing kernels with these basic constructs.
- Template meta-programming is the answer.
  - This is the same approach that Intel TBB and Thrust take.
  - Has the effect of requiring that Tpetra objects be templated on Node type.
- Node provides generic parallel constructs, user fills in the rest:

| | |
|---|---|
| ```template <class WDP>```<br>```void Node::parallel_for(```<br>  ```int beg, int end, WDP workdata);``` | ```template <class WDP>```<br>```WDP::ReductionType Node::parallel_reduce(```<br>  ```int beg, int end, WDP workdata);``` |
| Work-data pair (WDP) struct provides:<br>• loop body via `WDP::execute(i)` | Work-data pair (WDP) struct provides:<br>• reduction type `WDP::ReductionType`<br>• element generation via `WDP::generate(i)`<br>• reduction via `WDP::reduce(x,y)` |

Sandia National Laboratories

# Example Kernels: axpy() and dot()

```cpp
template <class WDP>
void
Node::parallel_for(int beg, int end,
                   WDP workdata    );
```

```cpp
template <class WDP>
WDP::ReductionType
Node::parallel_reduce(int beg, int end,
                      WDP workdata    );
```

```cpp
template <class T>
struct AxpyOp {
  const T * x;
  T * y;
  T alpha, beta;
  void execute(int i)
  { y[i] = alpha*x[i] + beta*y[i]; }
};
```

```cpp
template <class T>
struct DotOp {
  typedef T ReductionType;
  const T * x, * y;
  T identity()      { return (T)0;      }
  T generate(int i)  { return x[i]*y[i]; }
  T reduce(T x, T y) { return x + y;      }
};
```

```cpp
AxpyOp<double> op;
op.x = ...;  op.alpha = ...;
op.y = ...;  op.beta  = ...;
node.parallel_for< AxpyOp<double> >
            (0, length, op);
```

```cpp
DotOp<float> op;
op.x = ...;  op.y = ...;
float dot;
dot = node.parallel_reduce< DotOp<float> >
                (0, length, op);
```

Sandia National Laboratories

# *Hybrid CPU/GPU Computing*

# Hybrid Timings (Tpetra)

- Tests of a simple iterations:

  - power method: one sparse mat-vec, two vector operations

  - conjugate gradient: one sparse mat-vec, five vector operations

- DNVS/x104 from UF Sparse Matrix Collection (100K rows, 9M entries)

- NCCS/ORNL Lens node includes:

  - one NVIDIA Tesla C1060

  - one NVIDIA 8800 GTX

  - Four AMD quad-core CPUs

- Results are very tentative!

  - suboptimal GPU traffic

  - bad format/kernel for GPU

  - bad data placement for threads

| Node | PM (mflop/s) | CG (mflop/s) |
|---|---|---|
| Single thread | 140 | 614 |
| 8800 GPU | 1,172 | 1,222 |
| Tesla GPU | 1,475 | 1,531 |
| Tesla + 8800 | 981 | 1,025 |
| 16 threads | 816 | 1,376 |
| **1 node** 15 threads + Tesla | 867 | 1,731 |
| **2 nodes** 15 threads + Tesla | 1,677 | 2,102 |

Sandia National Laboratories

# New Core Linear Algebra Needs

# Advanced Modeling and Simulation Capabilities: Stability, Uncertainty and Optimization

- **Promise: 10-1000 times increase in parallelism (or more).**

**SPDEs:**

**Transient Optimization:**

Lower Block Bi-diagonal

$t_0$

$t_n$

Block Tri-diagonal

$t_0$

$t_n$

- **Pre-requisite: High-fidelity "forward" solve:**
  - **Computing families of solutions to similar problems.**
  - **Differences in results must be meaningful.**

■ **- Size of a single forward problem**

Sandia National Laboratories

# Advanced Capabilities: Readiness and Importance

| Modeling Area | Sufficient Fidelity? | Other concerns | Advanced capabilities priority |
|---|---|---|---|
| Seismic<br>*S. Collis, C. Ober* | Yes. | None as big. | Top. |
| Shock & Multiphysics (Alegra)<br>*A. Robinson, C. Ober* | Yes, but some concerns. | Constitutive models, material responses maturity. | Secondary now. Non-intrusive most attractive. |
| Multiphysics (Charon)<br><br>*J. Shadid* | Reacting flow w/ simple transport, device w/ drift diffusion, … | Higher fidelity, more accurate multiphysics. | Emerging, not top. |
| Solid mechanics<br><br>*K. Pierson* | Yes, but… | Better contact. Better timestepping. Failure modeling. | Not high for now. |

# Advanced Capabilities:
## Other issues

- Non-intrusive algorithms (e.g., Dakota):
  - Task level parallel:
    - A true peta/exa scale problem?
    - Needs a cluster of 1000 tera/peta scale nodes.
- Embedded/intrusive algorithms (e.g., Trilinos):
  - Cost of code refactoring:
    - Non-linear application becomes "subroutine".
    - Disruptive, pervasive design changes.
- Forward problem fidelity:
  - Not uniformly available.
  - Smoothness issues.
  - Material responses.

Sandia National Laboratories

# Advanced Capabilities:
## Derived Requirements

- Large-scale problem presents collections of related subproblems with forward problem sizes.

- Linear Solvers:  $Ax = b \rightarrow AX = B, \; Ax^i = b^i, \; A^i x^i = b^i$

  - Krylov methods for multiple RHS, related systems.

- Preconditioners:

  $$A^i = A_0 + \Delta A^i$$

  - Preconditioners for related systems.

- Data structures/communication:  $pattern(A^i) = pattern(A^j)$

  - Substantial graph data reuse.

# Summary

- App targets will change:
  - Advanced modeling and simulation: Gives a better answer.
  - Kernel set changes.
- Resilience requires an integrated strategy:
  - Most effort at the system/runtime level.
  - C/R (with localization) will continue at the app level.
  - Resilient algorithms will mitigate soft error impact.
- Building the next generation of parallel applications requires enabling domain scientists:
  - Write sophisticated methods.
  - Do so with serial fragments.
  - Fragments hoisted into scalable, resilient fragment.

Sandia National Laboratories

# Quiz (True or False)

1. MPI-only has the best parallel performance.
2. Future parallel applications will not have MPI_Init().
3. All future programmers will need to write parallel code.
4. Use of "markup", e.g., OpenMP pragmas, is the least intrusive approach to parallelizing a code.
5. DRY is not possible across CPUs and GPUs
6. GPUs are a harbinger of CPU things to come.
7. Checkpoint/Restart will be sufficient for scalable resilience.
8. Resilience will be built into algorithms.
9. MPI-only and MPI+X can coexist in the same application.
10. Kernels will be different in the future.